

Plastid genome assembly

Tomáš Fér

Dept. of Botany, Charles University, Prague

June 2025

Basic approaches

- mapping to reference genes/introns/spacers (HybPhyloMaker)
 - a script available to prepare HPM-compatible reference from GenBank files
- mapping to full-genome reference (separate script using HPM output)
- *de novo* assembly/annotation – many pipelines
 - FastPlast
 - GetOrganelle
 - NOVOplasty
 - ORG.asm
 - ...

Reference mapping

- prepare a reference from GenBank file using ‘`HybPhyloMakerOf_createPlastomeRef.sh`’

```
gene          complement(101946..102413)
/gene="rps7"
/locus_tag="LK123_pgp023"
/db_xref="GeneID:68666775"
complement(101946..102413)
/gene="rps7"
/locus_tag="LK123_pgp023"
/codon_start=1
/transl_table=11
/product="ribosomal protein S7"
/protein_id="YP\_010219699.1"
/db_xref="GeneID:68666775"
/translation="MSRRGTAEKTAKS DPIYRNRLVNMLVNRLKHGKKSLAYQIY
RAMKKIQQQKTETNPLSVLRQAIRGVTPDIAVKARRVSGSTHQVPIEIGSTQGKALAIR
WLLGASRKPRGRNMAFKLSELVDAAKGSGDAIRKKEETHRMAEANRAFAHFR"
105205..105276
/gene="trnV-GAC"
/locus_tag="LK123_pgt025"
/db_xref="GeneID:68666776"
105205..105276
/gene="trnV-GAC"
/locus_tag="LK123_pgt025"
/product="tRNA-Val"
/db_xref="GeneID:68666776"
105504..106994
/gene="rrn16"
/locus_tag="LK123_pgr008"
/db_xref="GeneID:68666777"
105504..106994
/gene="rrn16"
/locus_tag="LK123_pgr008"
/product="16S ribosomal RNA"
/db_xref="GeneID:68666777"
```

```
>109_109_trnMxCAU_tRNA
acctacttaactcagcggttagagtattgctttcatacggc
gggagtcattggttcaaatccaatagtaggtt
>110_110_trnMxCAU-atpE_non
gaacttatttagataccgcagtcaatggtatctaataagttt
ttatacacacattgattttagtaatattttttgtatcttt
>111_111_atpE_CDS
ttatgaaatggcattgatagcctctactcgtgtcctagccc
gtcggagagctagattgcctcaattgtttgtcttttcct
tcagctttctcaaaggcagctccgctattcaagagttt
>230_230_4.5S_rRNA
gaaggtcacggcgagacgagccgttattcattacgataggt
gtcaagtggaaagtgcagtgtatgcagctgaggcatcct
aacagaccggtagacttgaac
```

- use the reference within HybPhyloMaker

Mapping to full plastome reference

- use sequence of a full plastome as a reference in HybPhyloMaker
(one IR should be removed)
- check the lengths of LSC, IRa and SSC in the paper associated with the published chloroplast genome and extract only this sequence

De novo assembly

- extract chloroplast derived reads from whole genome or enriched dataset
- produce the assembly, ideally fully circular
- identify LSC, SSC and IRs
- annotate the assembly
 - DOGMA (<https://dogma.ccbb.utexas.edu/>) – discontinued...
 - GeSeq (<https://chlorobox.mpimp-golm.mpg.de/geseq.html>)
 - Plastid Genome Annotator (PGA) (<https://github.com/quxiaojian/PGA>)

Fast-Plast



- <https://github.com/mrmckain/Fast-Plast>
- read trimming (**Trimmomatic**)
- read extraction – mapping to reference (**Bowtie 2**)
- de novo assembly using de Bruijn graphs (**SPAdes**)
- iterative seed-based assembly to close gaps of contigs with low coverage (**afin**)
- check for gene content (**BLAST+**)
- identification of quadripartite structure and proper ordering of LSC, SSC and IRs
- coverage analysis (**Jellyfish 2**)

Fast-Plast



- input files – gzipped FASTQ files (PE or SE)
- specify plant order for plastome reads retrieval
(FastPlast includes more than 1,000 plastomes
to create Bowtie 2 index)
- --name – prefix to all files

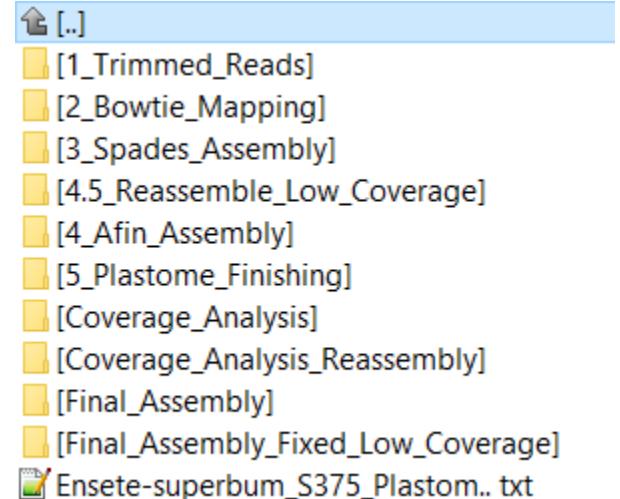
```
perl fast-plast.pl \
-1 genus-species_code_R1.fastq.gz \
-2 genus-species_code_R2.fastq.gz \
--name genus-species_code \
--bowtie_index Zingiberales \
--coverage_analysis
```



Fast-Plast results

- input files – gzipped FASTQ files (PE or SE)
- specify plant order for plastome reads retrieval
(FastPlast includes more than 1,000 plastomes to create Bowtie 2 index)
- --name – prefix to all files

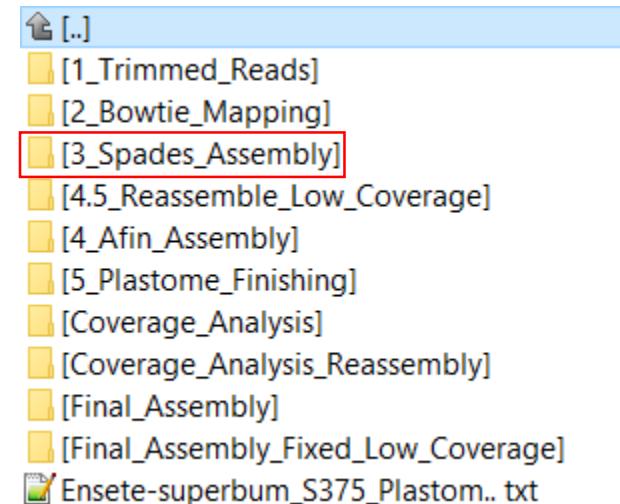
```
perl fast-plast.pl \
    -1 genus-species_code_R1.fastq.gz \
    -2 genus-species_code_R2.fastq.gz \
    --name genus-species_code \
    --bowtie_index Zingiberales \
    --coverage_analysis
```



Fast-Plast assembly



- SPAdes (<https://cab.spbu.ru/software/spades/>)
- de Bruijn graph assembler
- Prjibelski A, Antipov D, Meleshko D, Lapidus A, Korobeynikov A. 2020. [Using SPAdes de novo assembler. Current Protocols in Bioinformatics](#), 70, e102. doi: 10.1002/cpbi.102



de Bruijn graphs

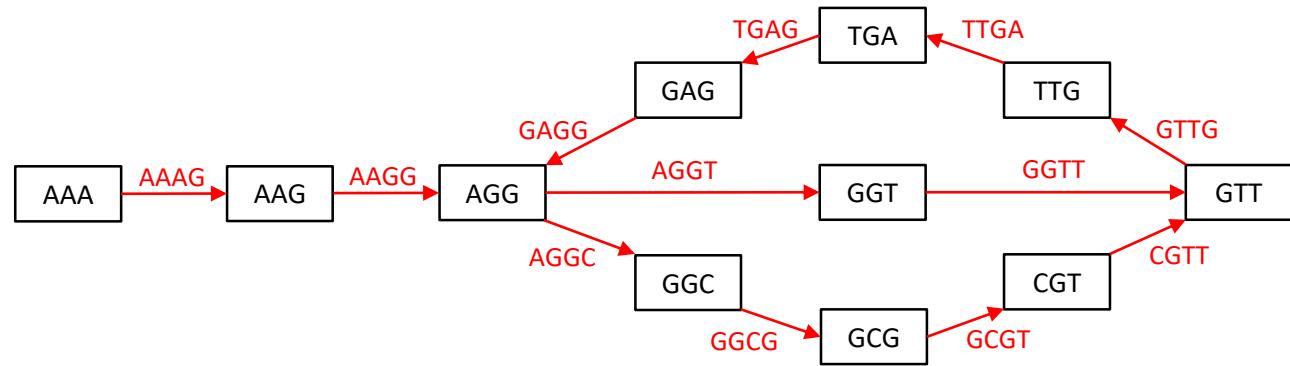
- network made up of nodes and edges (directed multigraph)
- these comes from the overlaps between k-mers

k-mers – a (DNA) molecule of the length k

- every possible $(k-1)$ -mer is assigned to a node
- edges are all possible k-mers
- connect nodes by a directed edge if there is a k-mer whose prefix (i.e., all position except the last one) is the former node
- suffix (i.e., all position except the first one) is the latter node
- Eulerian cycle in the graph (Eulerian walk) – visits each edge exactly ones

k -mer = 4

AAAGGCCTTGAGGTT
AAAG
AAGG
AGGC
GGCG
GCGT
CGTT
GTTG
TTGA
TGAG
GAGG
AGGT
GGCT
CGTT
GTTG



single or multiple Eulerian walk possible? Why?

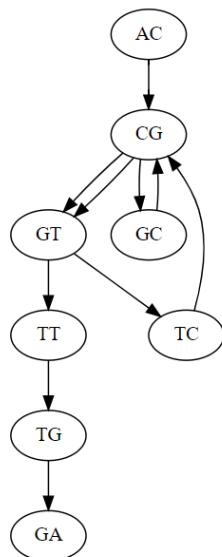
de Bruijn graphs

- requirements for straightforward graph
 - all k -mers present in the genome sequenced (gaps in sequencing lead to fragmented graphs)
 - all k -mers are error-free (error correction possible)
 - each k -mer appears at most once in the genome (different coverage requires normalization)
 - genome consists of a single circular chromosome
- play with k -mers and graphs using this Jupyter Notebook by B. Langmead
<https://colab.research.google.com/drive/1pQu9tJZ9RNpk8AaL2ThEYXoI3lu7Rw34>
- experiment with different k -mer settings

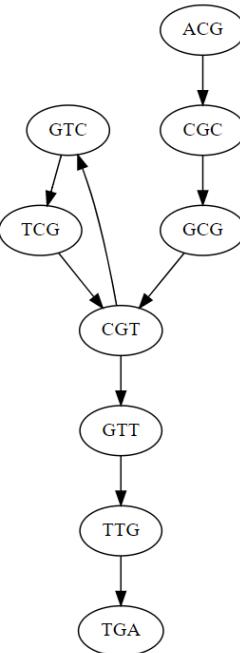
de Bruijn graphs

ACGCGTCGTTGA

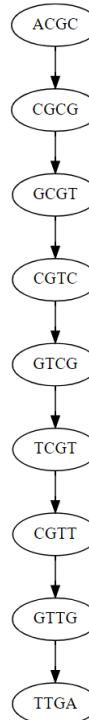
$k=3$



$k=4$



$k=5$



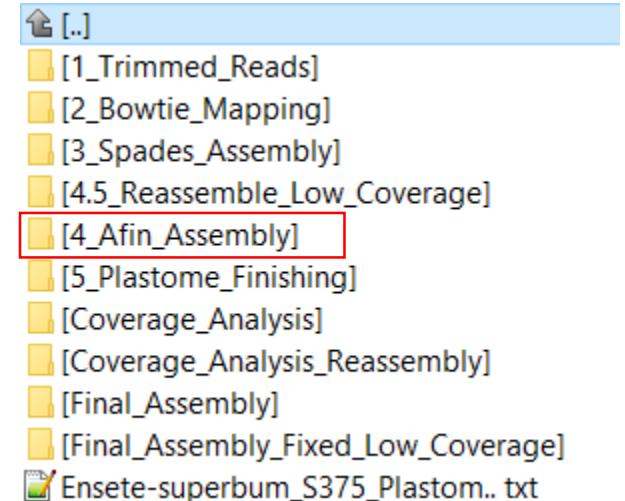
All graphs with Eulerian path

- all nodes (except first and last) are balanced (i.e., # incoming edges = # outgoing edges)
- starting and ending nodes are semibalanced



Fast-Plast assembly

- SPAdes (<https://cab.spbu.ru/software/spades/>)
 - de Bruijn graph assembler
- afin assembly on SPAdes contigs
 - iterative seed-based assembly – contig fusion and extension using the trimmed read set
 - for closing gaps
 - <https://github.com/mrmckain/Fast-Plast/tree/master/afin>
- scaffolding with SSPACE (https://github.com/nsoranzo/sspace_basic)
 - if more than one contig found after afin
 - contig extension/scaffolding using PE reads
 - Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W. 2011. [Scaffolding pre-assembled contigs using SSPACE](#). Bioinformatics, 27, 578-579





Fast-Plast finishing

- identification of genes in the assembly ([BLAST](#))
 - gene composition of the assembly
- identification/orientation of LSC, SSC and IRs
 - full sequences
 - sequences split into 4 pieces

[..]
[1_Trimmed_Reads]
[2_Bowtie_Mapping]
[3_Spades_Assembly]
[4.5_Reassemble_Low_Coverage]
[4_Afin_Assembly]
[5_Plasmome_Finishing]
[Coverage_Analysis]
[Coverage_Analysis_Reassembly]
[Final_Assembly]
[Final_Assembly_Fixed_Low_Coverage]
Ensete-superbum_S375_Plasmom.. txt



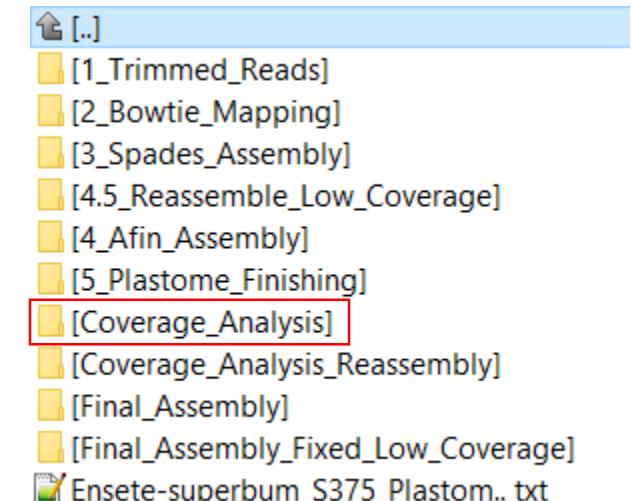
Fast-Plast coverage analysis

- k-mer coverage (Jellyfish 2)
- sudden coverage ‘doubling’ identifies IR boundary

79096	ACTTTACTCCTTTTTTACATT	79095	110
79097	CTTACTCCTTTTTTACATTT	79096	115
79098	TTTACTCCTTTTTTACATTTT	79097	114
79099	TTACTCCTTTTTTACATTTT	79098	114
79100	TACTCCTTTTTTACATTTTT	79099	112
79101	ACTCCTTTTTTACATTTTTA	79100	114
79102	CTCCTTTTTTACATTTTTAT	79101	114
79103	TCCTTTTTTACATTTTTATT	79102	112
79104	CCTTTTTTTTACATTTTTATT	79103	112
79105	CTTTTTTTTACATTTTTATT	79104	118
79106	TTTTTTTTTACATTTTTATTTC	79105	196
79107	TTTTTTTTTACATTTTTATTTC	79106	199
79108	TTTTTTTACATTTTTATTTC	79107	201
79109	TTTTTTTACATTTTTATTTC	79108	207
79110	TTTTTTACATTTTTATTTC	79109	208
79111	TTTTTACATTTTTATTTC	79110	205
79112	TTTTACATTTTTATTTC	79111	205
79113	TTTACATTTTTATTTC	79112	199
79114	TTACATTTTTATTTC	79113	198
79115	TACATTTTTATTTC	79114	204
79116	ACATTTTTATTTC	79115	209
79117	CATTTTTATTTC	79116	212
79118	ATTTTTATTTC	79117	209
79119	TTTTTATTTC	79118	207

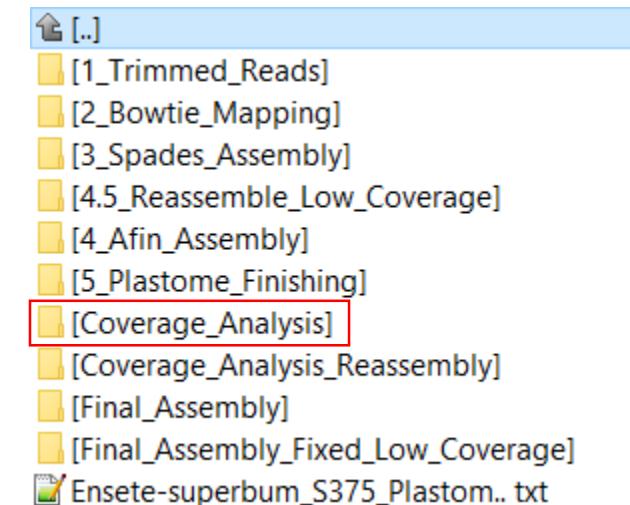
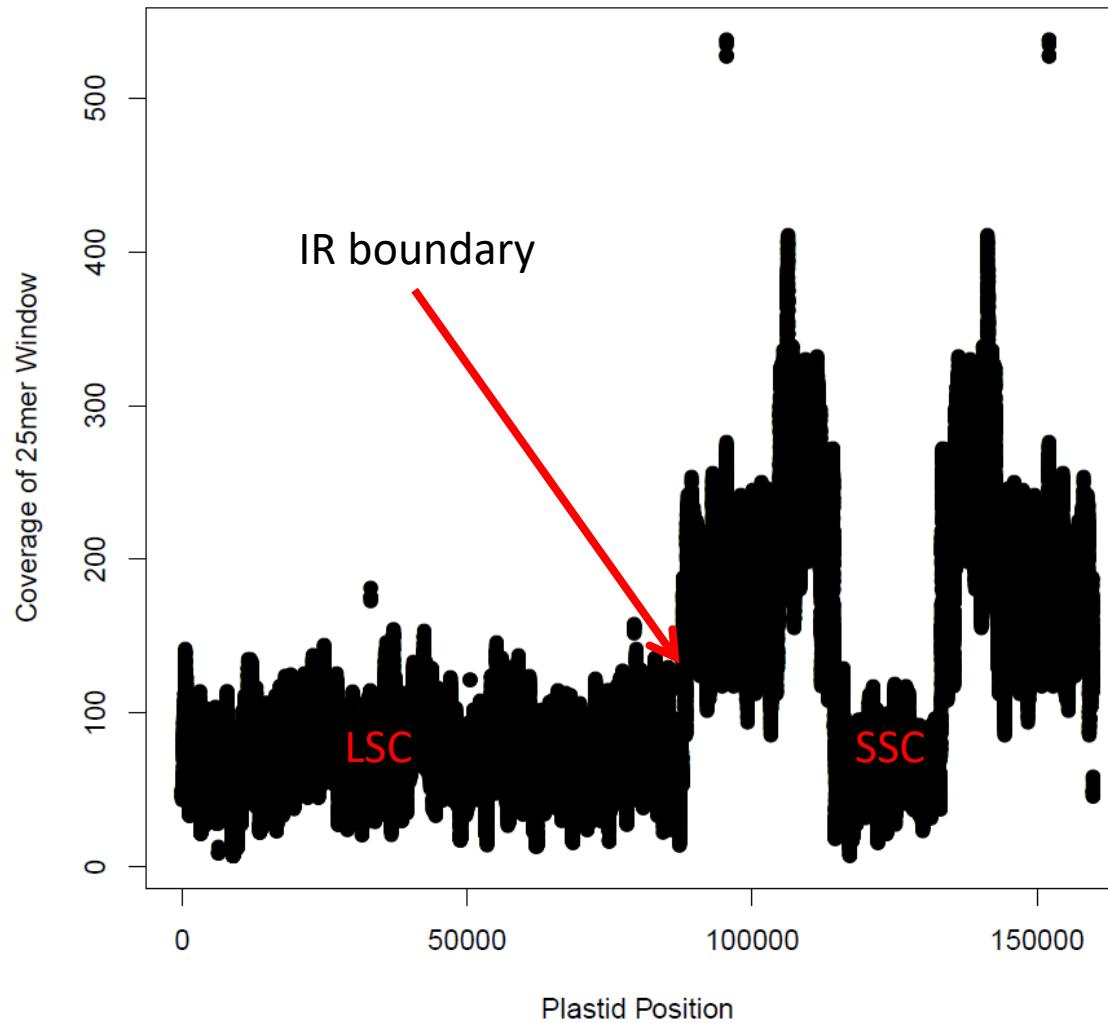
end of LSC

beginning of IR



Marçais G & Kingsford C. 2011. [A fast, lock-free approach for efficient parallel counting of occurrences of k-mers](#). *Bioinformatics*, 27, 764-770

Fast-Plast coverage analysis





Fast-Plast reassembly

- if regions with low coverage were identified

- low coverage regions removed
- contig broken into pieces
- reassembly from afin step
- coverage analysis of reassembly

[..]
[1_Trimmed_Reads]
[2_Bowtie_Mapping]
[3_Spades_Assembly]
[4.5_Reassemble_Low_Coverage]
[4_Afin_Assembly]
[5_Plasmome_Finishing]
[Coverage_Analysis]
[Coverage_Analysis_Reassembly]
[Final_Assembly]
[Final_Assembly_Fixed_Low_Coverage]
Ensete-superbum_S375_Plastom.. txt



Fast-Plast final results

Sample: Ensete-superbum_S375

Fast-Plast Version: Fast-Plast v.1.2.8

Total Cleaned Pair-End Reads: 1095096

Total Cleaned Single End Reads: 11531

Total Concordantly Mapped Reads: 199582

Total Non-concordantly Mapped Reads: 1273

Total Chloroplast Genome Length: 159320

Large Single Copy Size: 79105

Inverted Repeat Size: 34607

Small Single Copy Size: 11001

Minimum Coverage Used for Verification: 31.9663802410244

VALUES BELOW FROM REASSEMBLED PLASTOME

Total Chloroplast Genome Length: 162818

Large Single Copy Size: 74743

Inverted Repeat Size: 38537

Small Single Copy Size: 11001

Average Large Single Copy Coverage: 125

Average Inverted Repeat Coverage: 330

Average Small Single Copy Coverage: 109

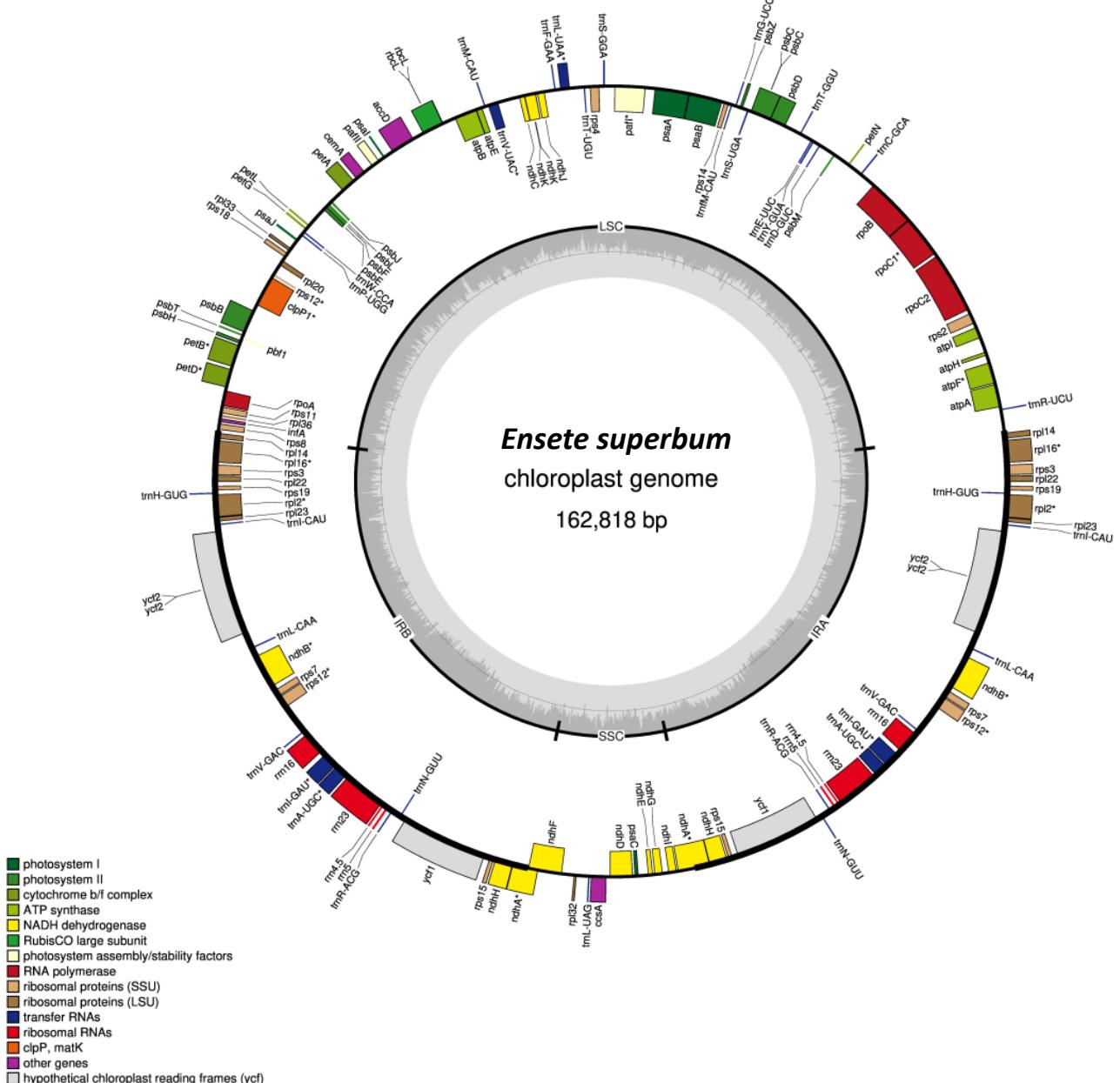
[..]
[1_Trimmed_Reads]
[2_Bowtie_Mapping]
[3_Spades_Assembly]
[4.5_Reassemble_Low_Coverage]
[4_Afin_Assembly]
[5_Plasmome_Finishing]
[Coverage_Analysis]
[Coverage_Analysis_Reassembly]
[Final_Assembly]
[Final_Assembly_Fixed_Low_Coverage]
Ensete-superbum_S375_Plastome.txt

Plastome annotation

- MPI-MP CHLOROBOX (<https://chlorobox.mpimp-golm.mpg.de/index.html>)
- GeSeq
 - upload FASTA file to annotate, reference possible
 - CDS, tRNA, rRNA
 - diverse tRNA annotators ([ARAGORN](#), [ARWEN](#), [tRNAscan-SE](#))
 - annotation support
 - Chloë (<https://chloe.plastid.org/>)
 - Mfannot (<https://megasun.bch.umontreal.ca/RNAweasel/>)
- results
 - output from primary annotation tools
 - annotation – GenBank, GFF3, GBSON
 - visualization – [OGDRAW](#)

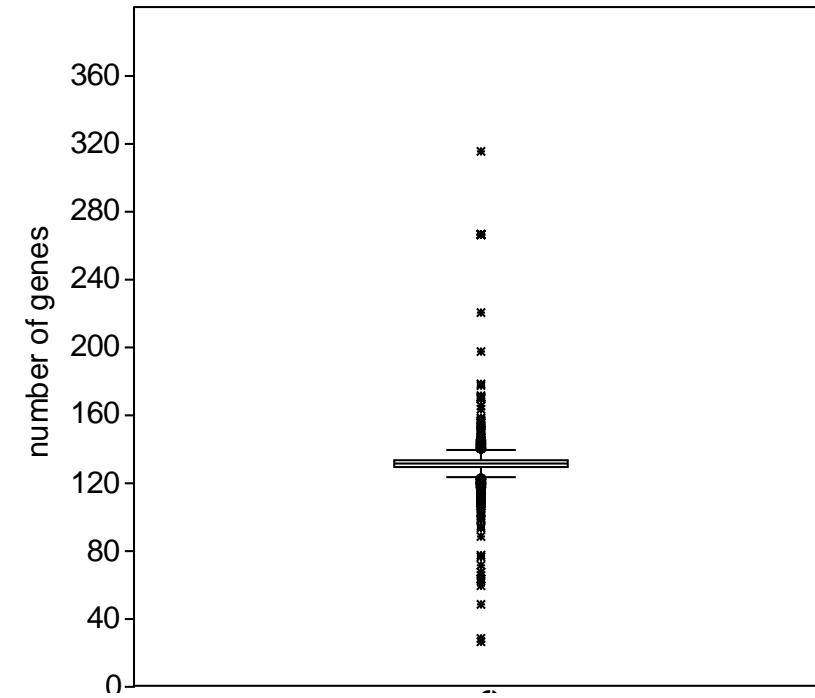
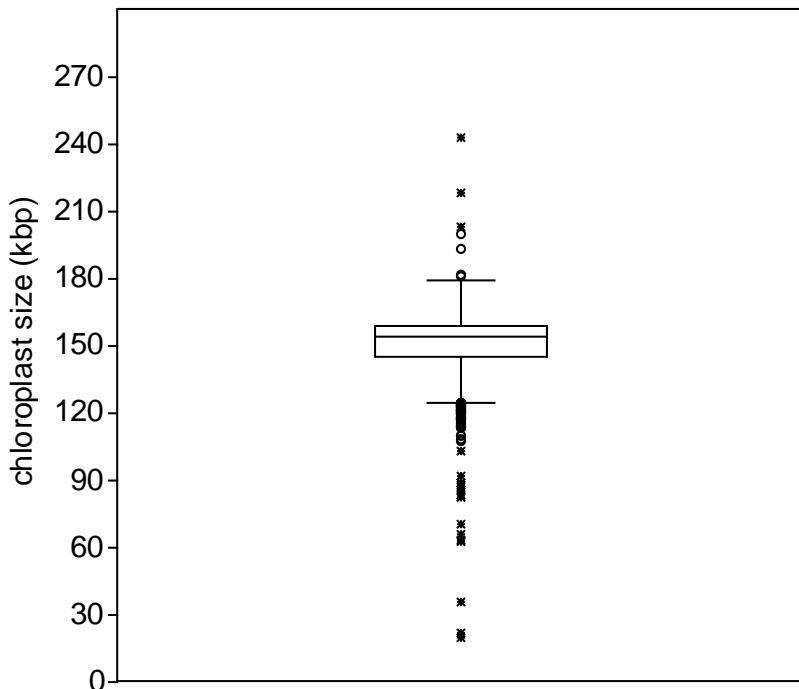
Plastome annotation

- OGDRAW - Draw Organelle Genome Maps
- <https://chlorobox.mpimp-golm.mpg.de/OGDraw.html>
- from GenBank file



Chloroplast size, number of genes

- cca 1,700 sequenced plastomes (land plants)
- size 150 kbp (19 – 243)
- 131 (26 – 315) genes: 84 proteins, 8 rRNA, 37 tRNA



Chloroplast size in lineages of land plants

